# Statistical Modelling III

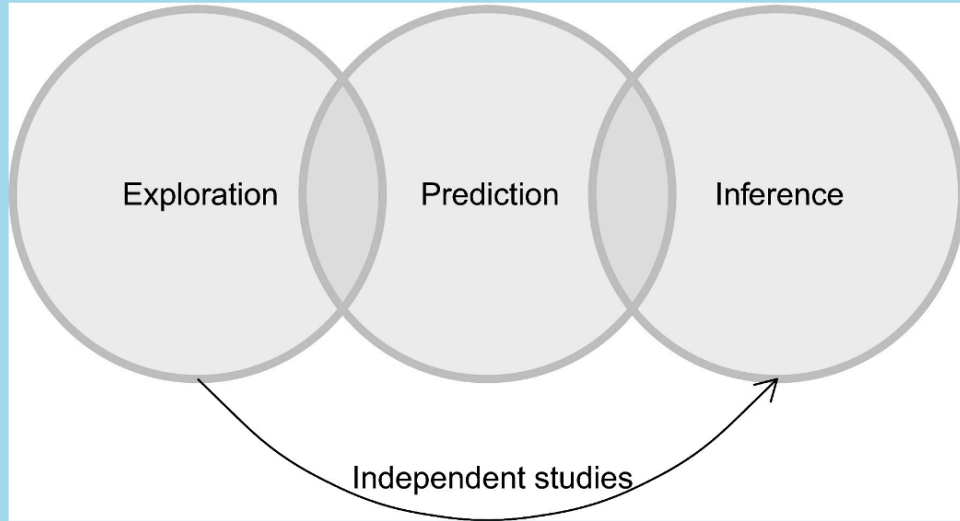HES 505 Fall 2024: Session 23

Carolyn Koehn

# Objectives

By the end of today you should be able to:

- Articulate three different reasons for modeling and how they link to assessments of fit

- Describe and implement several test statistics for assessing model fit

- Describe and implement several assessments of classification

- Describe and implement resampling techniques to estimate predictive performance

# The 3 Faces of Models

# Best Model for What?



from Tradennick et al. 2021

- **Exploration:** describe patterns in the data and generate hypotheses
- **Inference:** evaluate the strength of evidence for some statement about the process
- **Prediction:** forecast outcomes at unsampled locations based on covariates

# The Importance of Model Fit

- The general regression context:

$$\hat{y} = \mathbf{X}\hat{\beta}$$

- **Inference** is focused on robust estimates of $\hat{\beta}$ given the data we have

- **Prediction** is focused on accurate forecasts of $\hat{y}$ at locations where we have yet to collect the data

# Inference and Presence/Absence Data

- $\hat{\beta}$ is conditional on variables in the model **and** those not in the model

```r
1  nsamp <- 1000
2  df <- data.frame(x1 = rnorm(nsamp,0,1),
3                   x2 = rnorm(nsamp,0,1),
4                   x3 = rnorm(nsamp,0,1))
5
6  linpred <- 1 + 2*df$x1 -0.18*df$x2 -3.5*df$x3
7  y <- rbinom(nsamp, 1, plogis(linpred))
8  df <- cbind(df, y)
9
10 mod1 <- glm(y~x1 +x2, data=df, family="binomial")
11 mod2 <- glm(y~x1 +x2 + x3, data=df, family="binomial")
```

# Inference & Presence/Absence Data
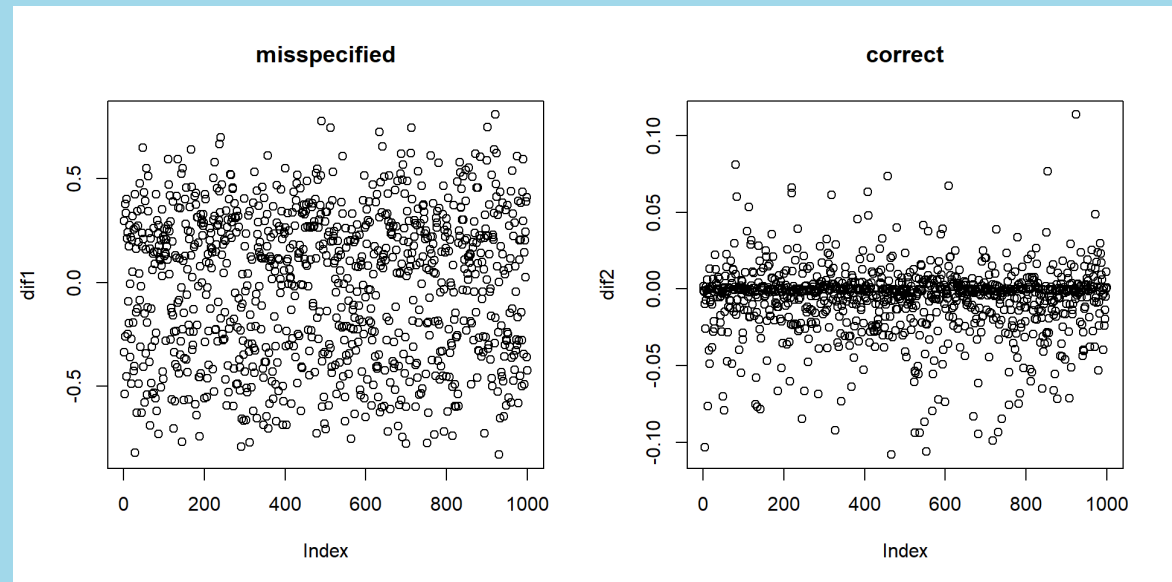
```
1  coef(mod1)
```

```
(Intercept)              x1
x2
   0.4460172    0.9041635
-0.1018929
```

```
1  coef(mod2)
```

```
(Intercept)              x1
x2              x3
   1.175262     2.298450
-0.307852    -3.774803
```
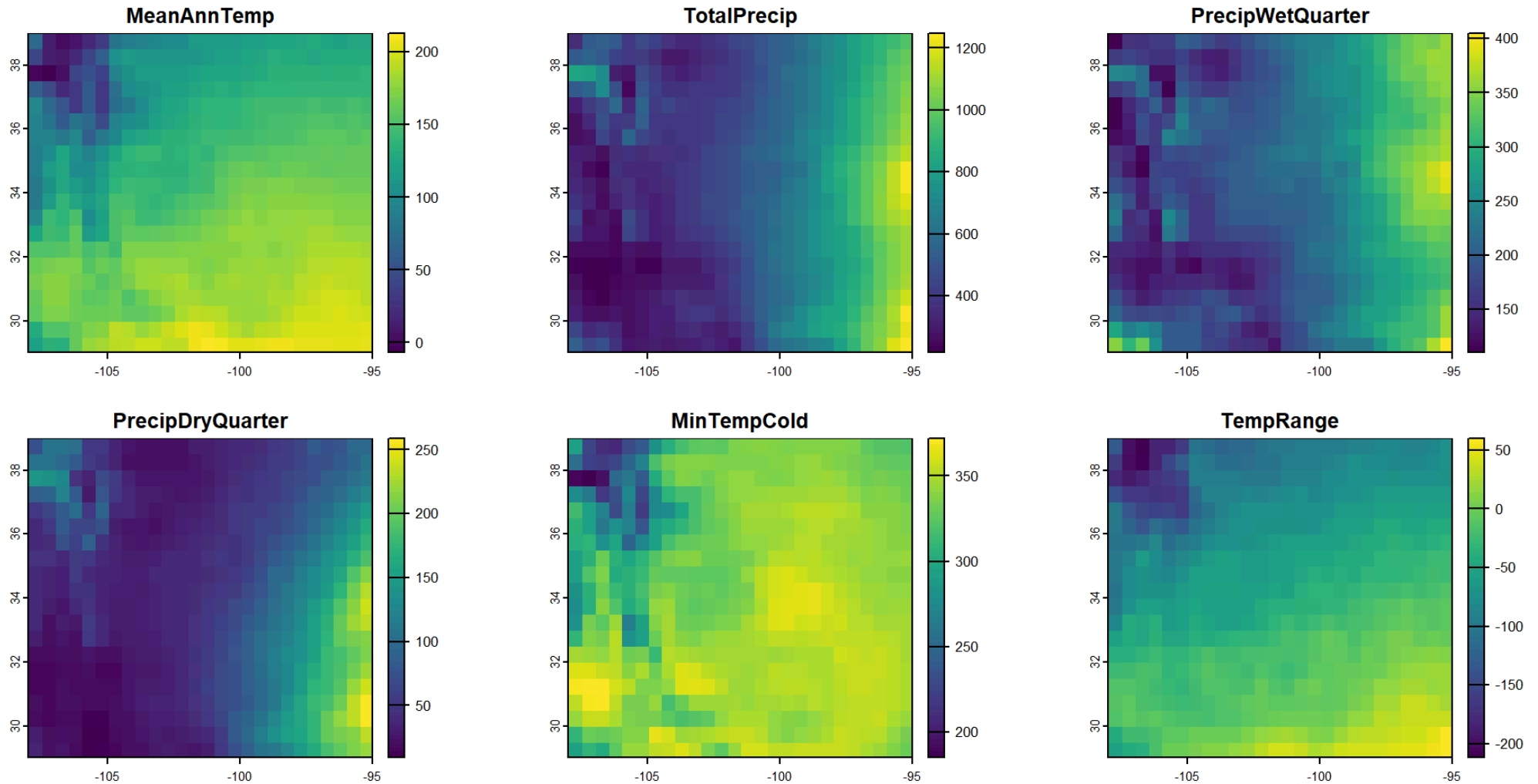
```
1  prd1 <- predict(mod1, df, "response")
2  dif1 <- plogis(linpred) - prd1
3  prd2 <- predict(mod2, df, "response")
4  dif2 <- plogis(linpred) - prd2
```



Inferring coefficient effects requires that your model fit the data well

# Assessing Model Fit

# Back to our simulated data

# Back to our simulated data

```r
1  base.path <- "/opt/data/data/presabsexample/" #sets the path to the root di
2
3  pres.abs <- st_read(paste0(base.path, "presenceabsence.shp"), quiet = TRUE)
4  pred.files <- list.files(base.path,pattern='grd$', full.names=TRUE) #get th
5
6  pred.stack <- rast(pred.files) #read into a RasterStack
7  names(pred.stack) <- c("MeanAnnTemp", "TotalPrecip", "PrecipWetQuarter", "P
8  plot(pred.stack)
9  pred.stack.scl <- scale(pred.stack)
10 pts.df <- terra::extract(pred.stack.scl, vect(pres.abs), df=TRUE)
11 pts.df <- cbind(pts.df, pres.abs$y)
12 colnames(pts.df)[8] <- "y"
```

# Using Test Statistics

- $R^2$ for linear regression:

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

$$SS_{res} = \sum_i (y_i - f_i)^2$$

$$SS_{tot} = \sum_i (y_i - \bar{y})^2$$

- Perfect prediction ($f_i = y_i$); $SS_{res} = 0$; and $R^2 = 1$

- Null prediction (Intercept only) ($f_i = \bar{y}$); $SS_{res} = SS_{tot}$; and $R^2 = 0$

- No direct way of implementing for logistic regression

# Pseudo- $R^2$

$$R_L^2 = \frac{D_{null} - D_{fitted}}{D_{null}}$$

- Cohen's Likelihood Ratio

- Deviance ($D$), the difference between the model and some hypothetical perfect model (lower is better)

- Challenge: Not monotonically related to $p$

- Challenge: How high is too high?

# Cohen's Likelihood Ratio

```
1  logistic.rich <- glm(y ~ MeanAnnTemp + PrecipWetQuarter + PrecipDryQuarter,
2                       family=binomial(link="logit"),
3                       data=pts.df[,2:8])
4  logistic.simple <- glm(y ~ MeanAnnTemp + TotalPrecip,
5                       family=binomial(link="logit"),
6                       data=pts.df[,2:8])
7
8  # Pseudo-R^2
9  with(logistic.rich,
10     null.deviance - deviance)/with(logistic.rich,
11                               null.deviance)
```

[1] 0.4495966

```
1  with(logistic.simple,
2      null.deviance - deviance)/with(logistic.simple,
3                               null.deviance)
```

[1] 0.4567641

# Pseudo- $R^2$

$$R^2_{CS} = 1 - \left( \frac{L_0}{L_M} \right)^{(2/n)}$$

$$= 1 - \exp^{2(ln(L_0)-ln(L_M))/n}$$

- Cox and Snell $R^2$
- Likelihood ($L$), the probability of observing the sample given an assumed distribution
- Challenge: Maximum value is less than 1 and changes with $n$
- Correction by Nagelkerke so that maximum is 1

# Cox and Snell $R^2$

```r
1  logistic.null <- glm(y ~ 1,
2                          family=binomial(link="logit"),
3                          data=pts.df[,2:8])
4
5  # Cox & Snell R^2 for logistic.rich
6  1 - exp(2*(logLik(logistic.null)[1] - logLik(logistic.rich)[1])/nobs(logist
```

[1] 0.4308873

```r
1  # Cox & Snell R^2 for logistic.simple
2  1 - exp(2*(logLik(logistic.null)[1] - logLik(logistic.simple)[1])/nobs(logi
```

[1] 0.4359785

# Using Test Statistics

- Based on the data used in the model (i.e., not prediction)
- Likelihood Ratio behaves most similarly to $R^2$
- Cox and Snell (and Nagelkerke) increases with more presences
- Ongoing debate over which is "best"
- **Don't defer to a single statistic**

# Assessing Predictive Ability

# Predictive Performance and Fit

- Predictive performance can be an estimate of fit

- Comparisons are often relative (better $\neq$ good)

- Theoretical and subsampling methods

# Theoretical Assessment of Predictive Performance

Hirotugu Akaike of AIC

- Information Criterion Methods
- Minimize the amount of information lost by using model to approximate true process
- Trade-off between fit and overfitting
- Can't know the true process (so comparisons are relative)

$$AIC = -2ln(\hat{L}) + 2k$$

# AIC Comparison

```
1 logistic.null$formula
```

y ~ 1

```
1 logistic.rich$formula
```

y ~ MeanAnnTemp + PrecipWetQuarter + PrecipDryQuarter

```
1 logistic.simple$formula
```
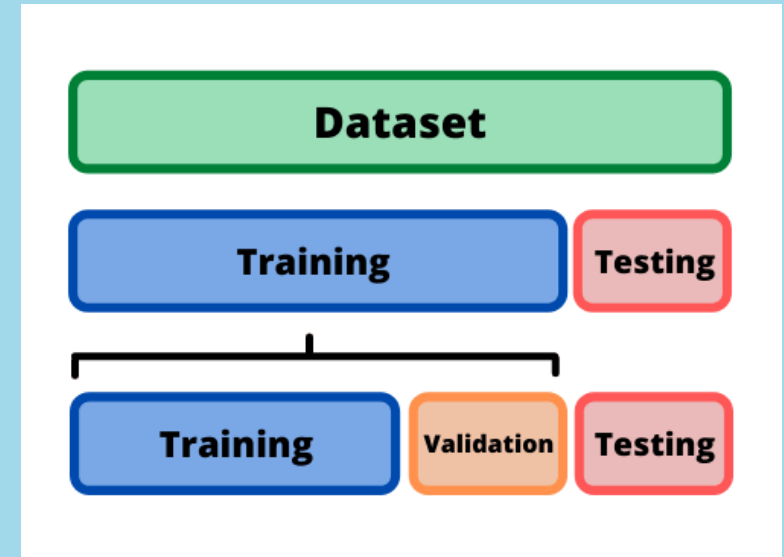
y ~ MeanAnnTemp + TotalPrecip

```
1 AIC(logistic.null, logistic.rich, logistic.simple)
```

```
                 df        AIC
logistic.null     1 127.37389
logistic.rich     4  77.00622
logistic.simple   3  74.10760
```

# Sub-sampling Methods

- Split data into *training* and *testing*

- Testing set needs to be large enough for results to be statistically meaningful

- Test set should be representative of the data as a whole

- Validation data used to tune parameters (not always)

# Subsampling your data with **caret**

```
1  pts.df$y <- factor(ifelse(pts.df$y == 1, "Yes", "No"),
2                       levels = c("Yes", "No"))
3  library(caret)
4  Train <- createDataPartition(pts.df$y, p=0.6, list=FALSE)
5
6  training <- pts.df[ Train, ]
7  testing <- pts.df[ -Train, ]
```

# Misclassification

- Confusion matrices compare actual values to predictions

- True Positive (TN) - This is correctly classified as the class if interest / target.

- True Negative (TN) - This is correctly classified as not a class of interest / target.

- False Positive (FP) - This is wrongly classified as the class of interest / target.

- False Negative (FN) - This is wrongly classified as not a class of interest / target.

|  | Actual Values | |
|---|---|---|
|  | Positive (1) | Negative (0) |
| Predicted Values Positive (1) | TP | FP |
| Predicted Values Negative (0) | FN | TN |

# Confusion Matrices in R

```r
train.log <- glm(y ~ .,
                 family="binomial"
                 data=training[,2:
predicted.log <- predict(train.log
                         newdata=t
                         type="res

pred <- factor(
  ifelse(predicted.log > 0.5,
                         "Yes",
                         "No"),
  levels = c("Yes", "No"))
```

```r
confusionMatrix(testing$y, pred)
```

```
Confusion Matrix and Statistics

          Reference
Prediction Yes No
       Yes   5  7
       No   26  1

               Accuracy : 0.1538
                 95% CI : (0.0586, 0.3053)
    No Information Rate : 0.7949
    P-Value [Acc > NIR] : 1.000000

                  Kappa : -0.3794

 Mcnemar's Test P-Value : 0.001728

            Sensitivity : 0.16129
            Specificity : 0.12500
         Pos Pred Value : 0.41667
         Neg Pred Value : 0.03704
             Prevalence : 0.79487
         Detection Rate : 0.12821
```

# Confusion Matrices

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Depends upon threshold!!

$$\text{Sensitivity (Recall)} = \frac{TP}{TP + FN}$$

$$\text{Specificity} = \frac{TN}{FP + TN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

# Confusion Matrices in R

```
1  library(tree)
2  tree.model <- tree(y ~ . , trainin
3  predict.tree <- predict(tree.model
```

```
1  confusionMatrix(testing$y, predict.tree)
```

```
Confusion Matrix and Statistics

          Reference
Prediction Yes No
      Yes   6  6
      No    3 24

               Accuracy : 0.7692
                 95% CI : (0.6067, 0.8887)
    No Information Rate : 0.7692
    P-Value [Acc > NIR] : 0.5882

                  Kappa : 0.4179

 Mcnemar's Test P-Value : 0.5050

            Sensitivity : 0.6667
            Specificity : 0.8000
         Pos Pred Value : 0.5000
         Neg Pred Value : 0.8889
             Prevalence : 0.2308
         Detection Rate : 0.1538
```

# Confusion Matrices in R

```r
1  library(randomForest, quietly = TR
2  class.model <- y ~ .
3  rf <- randomForest(class.model, da
4  predict.rf <- predict(rf, newdata=
```

```r
1  confusionMatrix(testing$y, predict.rf)
```
Confusion Matrix and Statistics

```
           Reference
Prediction Yes No
      Yes   7  5
      No    4 23

               Accuracy : 0.7692
                 95% CI : (0.6067, 0.8887)
    No Information Rate : 0.7179
    P-Value [Acc > NIR] : 0.3037

                  Kappa : 0.4455

 Mcnemar's Test P-Value : 1.0000

            Sensitivity : 0.6364
            Specificity : 0.8214
         Pos Pred Value : 0.5833
         Neg Pred Value : 0.8519
             Prevalence : 0.2821
         Detection Rate : 0.1795
```
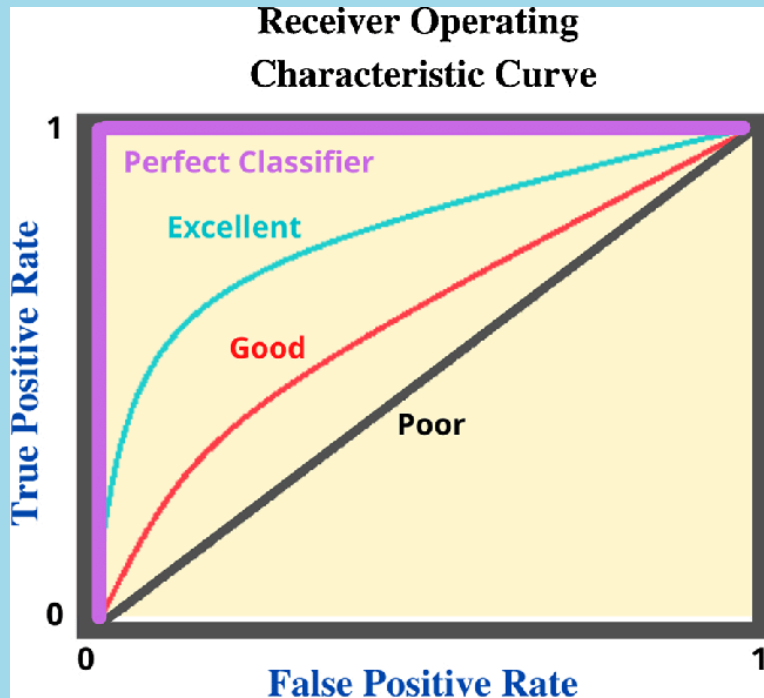
# Threshold-Free Methods



Receiver Operating Characteristic Curve

- Receiver Operating Characteristic Curves

- Illustrates discrimination of binary classifier as the threshold is varied

- Area Under the Curve (AUC) provides an estimate of classification ability

# Criticisms of ROC/AUC

- Treats false positives and false negatives equally

- Undervalues models that predict across smaller geographies

- Focus on *discrimination* and not *calibration*

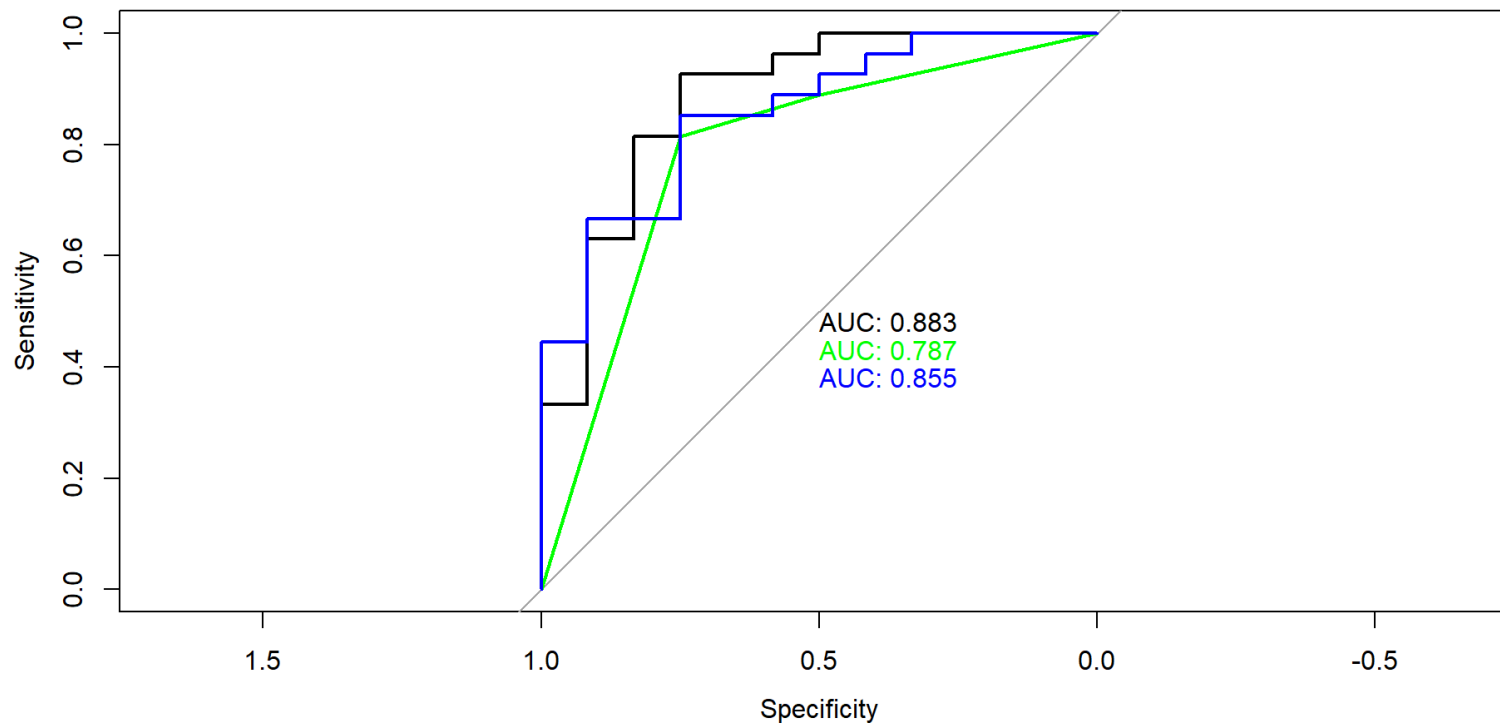- New methods for presence-only data

# ROC in R (using pROC)

- Generate predictions (note the difference for tree and rf)

```
 1  library(pROC, quietly = TRUE)
 2  train.log <- glm(y ~ .,
 3                    family="binomial",
 4                    data=training[,2:8])
 5
 6  predicted.log <- predict(train.log,
 7                           newdata=testing[,2:8],
 8                           type="response")
 9
10  predict.tree <- predict(tree.model, newdata=testing[,2:8], type="vector")[,
11
12  predict.rf <- predict(rf, newdata=testing[,2:8], type="prob")[,2]
```

# ROC in R (using pROC)

```
1  plot(roc(testing$y, predicted.log), print.auc=TRUE)
2
3  plot(roc(testing$y, predict.tree), print.auc=TRUE, print.auc.y = 0.45, col=
4
5  plot(roc(testing$y, predict.rf), print.auc=TRUE, print.auc.y = 0.4, col="bl
```



AUC: 0.883
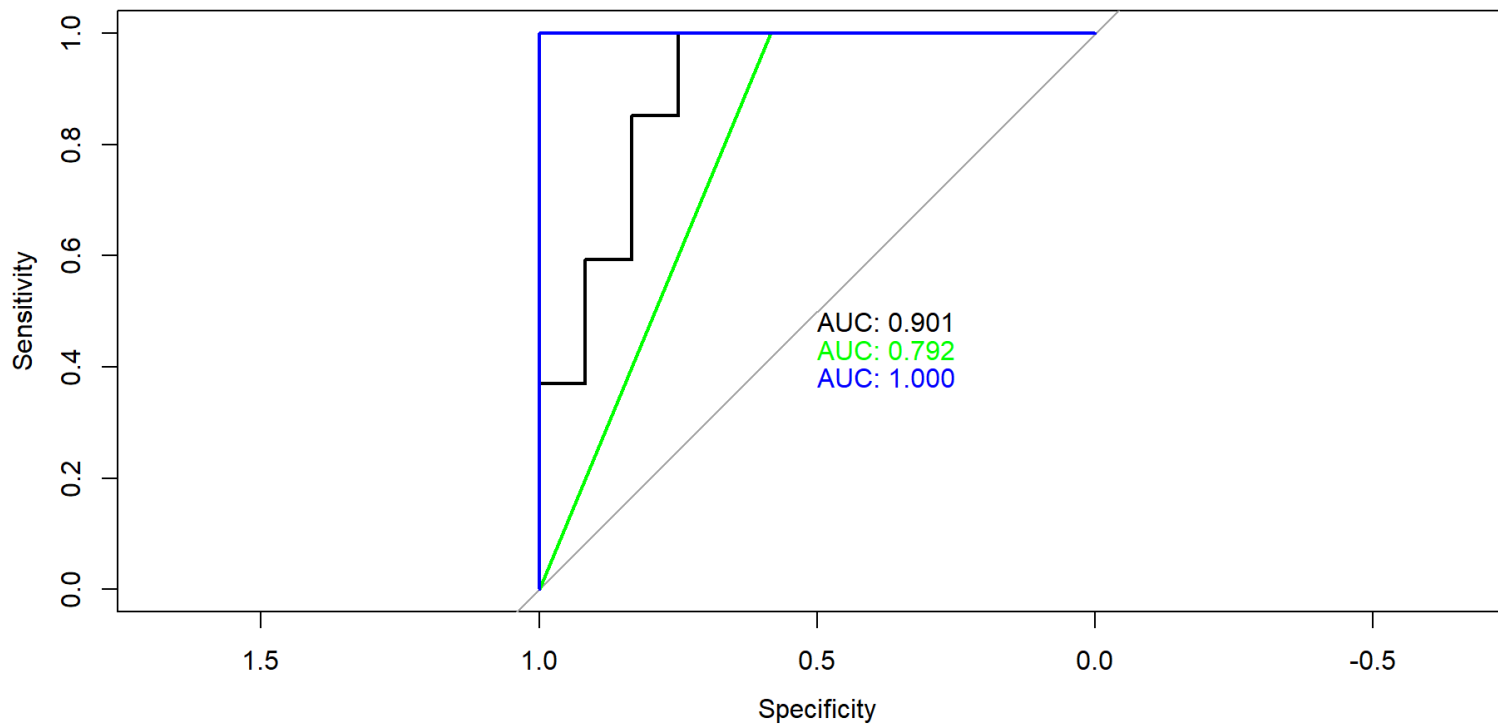AUC: 0.787
AUC: 0.855

# Cross-validation

- Often want to make sure that fit/accuracy not a function of partition choice

- Cross-validation allows resampling of data (multiple times)

- K-fold - Data are split into K datasets of ~ equal size, model fit to $(K-1)(\frac{n}{K})$ observations to predict held-out set

- Leave One Out (LOO) - Model fit to n-1 observations to predict the held out observation

# Crossvalidation in R using caret

```r
 1  fitControl <- trainControl(method = "repeatedcv",
 2                                  number = 10,
 3                                  repeats = 10,
 4                                  classProbs = TRUE,
 5                                  summaryFunction = twoClassSummary)
 6
 7  log.model <- train(y ~., data = pts.df[,2:8],
 8                  method = "glm",
 9                  trControl = fitControl)
10  pred.log <- predict(log.model, newdata = testing[,2:8], type="prob")[,2]
11
12  tree.model <- train(y ~., data = pts.df[,2:8],
13                  method = "rpart",
14                  trControl = fitControl)
15
16  pred.tree <- predict(tree.model, newdata=testing[,2:8], type="prob")[,2]
17
18  rf.model <- train(y ~., data = pts.df[,2:8],
```
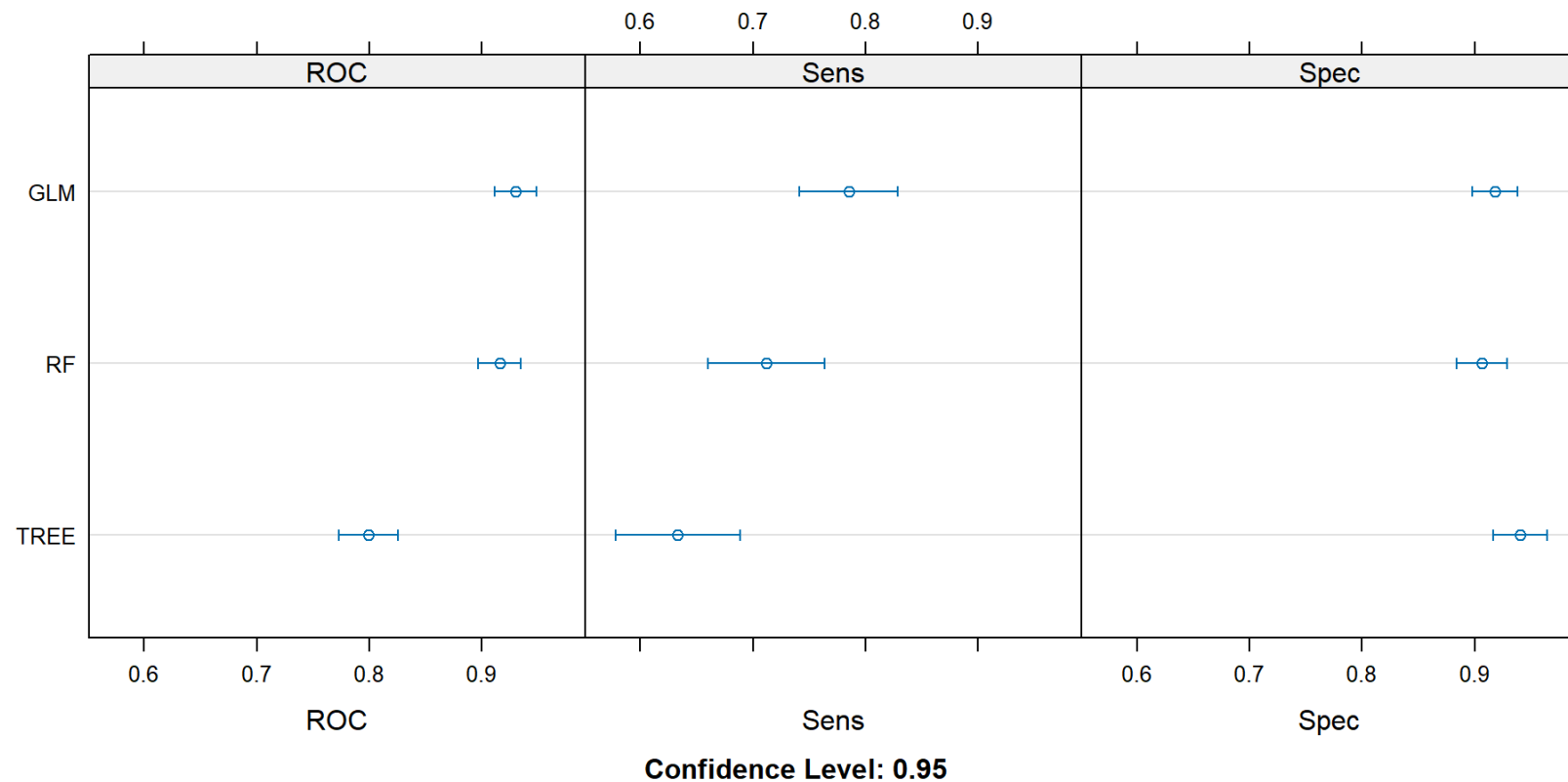
# Crossvalidation in **R** using **caret**

```r
1  plot(roc(testing$y, pred.log), print.auc=TRUE)
2
3  plot(roc(testing$y, pred.tree), print.auc=TRUE, print.auc.y = 0.45, col="gr
4
5  plot(roc(testing$y, pred.rf), print.auc=TRUE, print.auc.y = 0.4, col="blue"
```

# Crossvalidation in R using caret

```r
1  resamps <- resamples(list(GLM = log.model,
2                            TREE = tree.model,
3                            RF = rf.model))
4  dotplot(resamps)
```



Confidence Level: 0.95

# Spatial predictions

```
1  best.rf <- rf.model$finalModel
2  best.glm <- log.model$finalModel
3
4  rf.spatial <- terra::predict(pred.stack.scl, best.rf, type="prob")
5
6
7  glm.spatial <- terra::predict(pred.stack.scl, best.glm,type="response" )
```

# Spatial predictions